

TimeBlind: A Spatio-Temporal Compositionality Benchmark for Video LLMs

Baiqi Li^{1*} Kangyi Zhao² Ce Zhang¹ Chancharik Mitra³
Jean de Dieu Nyandwi³ Gedas Bertasius¹

¹University of North Carolina at Chapel Hill ²University of Pittsburgh ³Carnegie Mellon University
baiqili@cs.unc.edu

Abstract

Fine-grained spatio-temporal understanding is essential for video reasoning and embodied AI. Yet, while Multimodal Large Language Models (MLLMs) master static semantics, their grasp of temporal dynamics remains brittle. We present TimeBlind, a diagnostic benchmark for compositional spatio-temporal understanding. Inspired by cognitive science, TimeBlind categorizes fine-grained temporal understanding into three levels: recognizing atomic events, characterizing event properties, and reasoning about event interdependencies. Unlike benchmarks that conflate recognition with temporal reasoning, TimeBlind leverages a minimal-pairs paradigm: video pairs share identical static visual content but differ solely in temporal structure, utilizing complementary questions to neutralize language priors. Evaluating over 20 state-of-the-art MLLMs (e.g., GPT-5, Gemini 3 Pro) on 600 curated instances (2400 video-question pairs), reveals that the Instance Accuracy (correctly distinguishing both videos in a pair) of the best performing MLLM is only 48.2%, far below the human performance (98.2%). These results demonstrate that even frontier models lack temporal reasoning, positioning TimeBlind as a vital diagnostic tool for next-generation video understanding. We will release the data and code. Dataset and code are available at https://baiqi-li.github.io/timeblind_project/.

1. Introduction

Fine-grained spatio-temporal understanding is fundamental for long-horizon video reasoning [33, 43] and embodied AI [39, 50]. Beyond simply recognizing *what* is present, an intelligent system must infer *what changes*, *how it changes*, and *how multiple changes compose* into causal structures. While recent Multimodal Large Language Models (MLLMs) [4, 8, 9, 32, 35] demonstrate impressive performance on general benchmarks, their “sense of time” remains surprisingly brittle. As shown in Figure 1, even frontier mod-

*Corresponding author.

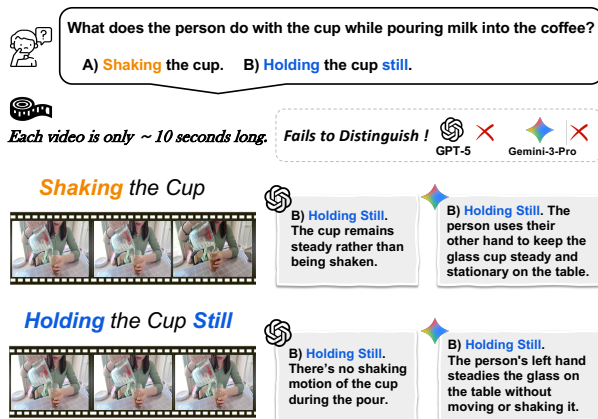


Figure 1. An example video pair that shares identical static visual content but differs solely in motion dynamics. The top video shows a person shaking a cup, while the bottom video shows them holding it still. Even the most advanced models like GPT-5 and Gemini 3 Pro fail to distinguish the actions in the video pair.

els (e.g., GPT-5 and Gemini 3 Pro) struggle to distinguish atomic actions (e.g., *shaking* vs. *holding*) in videos as short as 10 seconds. Moreover, models often misjudge relative dynamics (e.g., *accelerating* vs. *decelerating*), and frequently fail to resolve linguistic temporal connectives (e.g., “*as soon as*”). This discrepancy suggests that current benchmark scores [10, 13, 14, 24, 53, 56] severely overestimate genuine temporal understanding capabilities.

The root issue lies in evaluation design. Existing benchmarks [10, 24, 27, 53] rarely *isolate* temporal structure as the sole discriminative factor. As a result, many models exploit “static shortcuts”—correlating visual entities with answers without modeling time [19, 21]. Furthermore, language priors [23] allow models to guess answers based on textual plausibility. To truly *diagnose* understanding, an evaluation must hold the visual content constant and vary only the temporal dynamics.

To address this, we introduce **TimeBlind**, a diagnostic benchmark for compositional spatio-temporal understanding (Table 1). TimeBlind adopts a minimal-pairs design: each

Table 1. **Comparison with prior temporal video benchmarks.** We compare coverage of temporal reasoning categories and anti-shortcut design features. TimeBlind uniquely covers all 13 Allen temporal relations[†], includes causal and comparative reasoning, and employs both paired videos and complementary questions to eliminate static and linguistic shortcuts. The rightmost column shows the gap between human and best model performance, indicating benchmark difficulty.

BENCHMARK	SIZE	TEMPORAL REASONING COVERAGE					ANTI-SHORTCUT DESIGN		HUMAN-MODEL GAP
		EVENT TYPES	EVENT ATTRIBUTES	TEMPORAL TOPOLOGIES	CAUSAL	CROSS-EVENT COMPARISON	VIDEO PAIRS	COMPL. QUESTIONS	
MVBench [24]	4K	2	2	2	✗	✗	✗	✗	–
TOMATO [36]	1.4K	2	4	2	✗	✗	✗	✗	57.3%
Vinoground [54]	2K	1	1	2	✗	✗	✓	✗	40.0%
TempCompass [27]	4K	2	2	2	✗	✓	✓	✗	–
TimeBlind (Ours)	2.4K	2	6	13 [†]	✓	✓	✓	✓	50.0%

[†]Covers all 13 event relations in Allen’s Interval Algebra [1]: *before, after, meets, met-by, overlaps, overlapped-by, starts, started-by, finishes, finished-by, during, contains, equals*.

instance contains two videos with near-identical static visual content that differ *solely* in temporal structure. To eliminate language priors, we use complementary questions where the correct answer flips between paired videos. By removing static and linguistic shortcuts, TimeBlind forces models to rely exclusively on temporal evidence. Unlike large-scale noisy benchmarks, TimeBlind prioritizes high-fidelity diagnostic precision. Similar to Winoground [41], each instance serves as a rigorous test for a specific cognitive primitive, where high-quality annotations are prioritized over scale.

Crucially, unlike prior work [19, 54], TimeBlind evaluates logical composition, not just perception. Inspired by cognitive science, we extend the theory of *Image Compositionality* [18, 22] to the temporal domain, organizing the benchmark around a hierarchical taxonomy: (i) *Events*, addressing what happened (e.g., changes in object attributes or action understanding); (ii) *Event Attributes*, describing how events unfold (e.g., speed, magnitude of change); and (iii) *Structural Event Logic*, examining how multiple events are composed (e.g., temporal topology, causality, and cross-event comparison). Within this framework, we encompass a diverse set of 11 fine-grained categories in real-world scenarios. For example, regarding *Temporal Topology* in *Structural Event Logic*, we incorporate Allen’s Interval Algebra [1] by constructing videos that cover all 13 temporal event relations (e.g., *overlaps, meets, equals*) and designing questions that distinguish specific topology differences within each pair, going beyond simple sequencing relations (e.g., *before, after*) used in prior work.

We evaluate over 20 state-of-the-art MLLMs, including leading proprietary models (e.g., Gemini 3 Pro [35], GPT-5 [32]) as well as strong open-source models (e.g., Qwen3-VL [4], Molmo2 [9], PLM [8]), on 600 curated instances (totaling 2,400 video-question pairs). We find that even the top-performing model, Gemini-3 Pro, achieves only 48.2% Instance Accuracy (I-Acc), which requires correctly distinguishing both videos in a pair, falling far below human per-

formance (98.2%). Furthermore, while advanced models like Gemini-3 Pro and GPT-5 demonstrate reasonable proficiency in recognizing isolated events (achieving 49.2% and 58.3% I-Acc, respectively), their performance degrades sharply on event attributes—such as Speed (slowly vs. rapidly), and Force (forcefully vs. gently)—dropping to 36.7% and 32.3% I-Acc. Moreover, extensive ablation studies demonstrate that model performance remains poor even with increased input frames or test-time reasoning, as GPT-5 gains only 3.3% I-Acc, indicating models remain “time-blind” and fragile in temporal reasoning. We believe TimeBlind will serve as a vital diagnostic tool for evaluating and developing MLLMs capable of genuine temporal logic.

In summary, our main contributions are:

- **TimeBlind Benchmark.** A diagnostic, minimal-pairs benchmark that isolates temporal structure while minimizing static shortcuts and language priors.
- **Taxonomy of Temporal Compositionality.** A cognitive hierarchy—*Events, Attributes, and Structural Logic*—that guides systematic benchmark construction and evaluation.
- **Diagnostic Findings.** An evaluation of over 20 SOTA MLLMs revealing a significant gap between perceived and actual temporal reasoning capabilities.

2. Related Work

VideoQA Benchmarks. Early VideoQA datasets [15, 20, 34, 44, 46, 47, 53] focus on simple scenarios with short clips and limited question types. Recent benchmarks address comprehensive evaluation [10, 24, 29], complex reasoning [7, 13, 31, 38], long-form understanding [30, 37, 43, 45, 49, 56], and domain-specific settings [28, 33, 51]. However, most do not isolate temporal structure as the sole discriminative factor, allowing models to exploit static shortcuts—relying on object co-occurrence or language priors without genuinely modeling temporal dynamics [19, 21, 23].

Evaluating Spatio-Temporal Compositionality. To assess the spatio-temporal understanding capabilities of MLLMs,

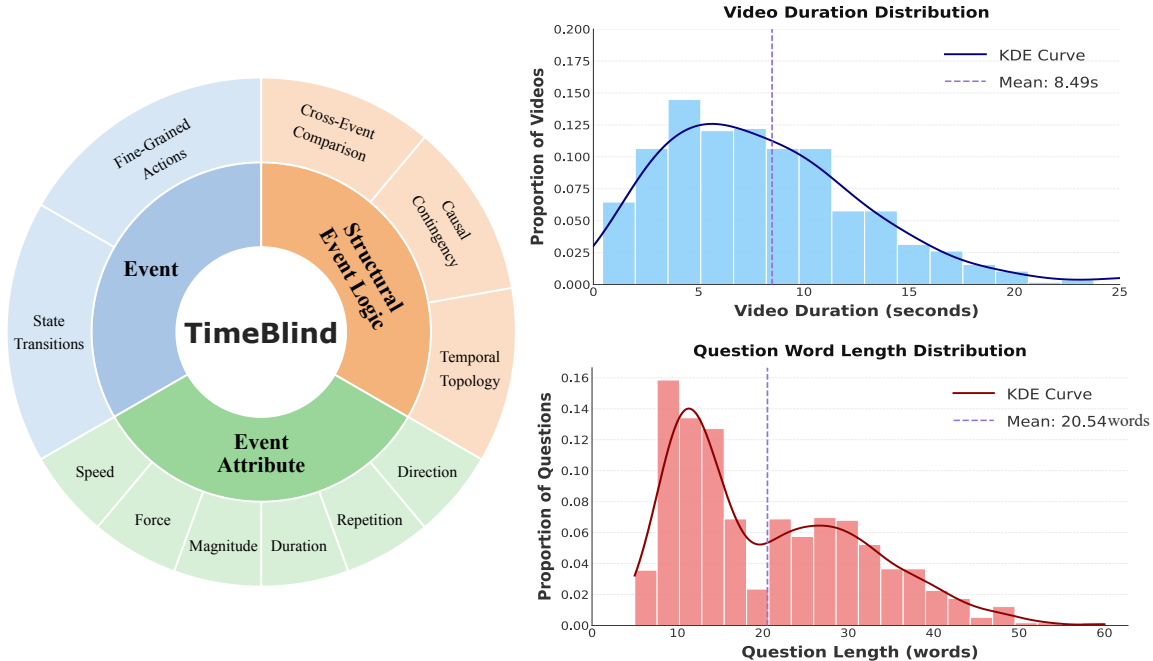


Figure 2. **TimeBlind Taxonomy and Statistics.** **Left:** We structure the evaluation into 11 fine-grained spatio-temporal compositional categories spanning three high-level aspects: Atomic Events (*what changes*), Parametric Event Attributes (*how it changes*), and Structural Event Logic (*how events compose*). **Top Right:** Distribution of video lengths across the benchmark, showing that most videos fall within the 0–15 seconds. **Bottom Right:** Distribution of question word counts, indicating that most questions are under 30 words. Overall, our benchmark features a structured taxonomy with diverse categories while maintaining short videos and concise questions.

several recent benchmarks have introduced diverse aspects of visual reasoning. In the image understanding domain, BLINK [11] reformats and groups classic vision problems into perception-centric multiple-choice questions and carefully removes language priors. In the video understanding domain, several works explicitly target temporal understanding by designing temporally challenging questions [5, 16, 36, 48]. Beyond single-video evaluation, a growing line of work adopts paired video–question protocols to more precisely diagnose temporal reasoning. TEMP-COMPASS [27] constructs paired videos by systematically manipulating an original video, such as reversing playback or altering temporal speed, thereby isolating temporal understanding from “static shortcuts.” VINOGROUND [54] further strengthens this paradigm by requiring models to answer identical questions over paired videos, where the correct answer is determined solely by temporal differences, effectively removing language priors. Follow-up works like GLIMPSE [57] and MVP [19] scale this approach to test physical and visual-centric reasoning. Different from prior works, TIMEBLIND achieves spatio-temporal compositionality through a carefully curated formal taxonomy. Drawing inspiration from cognitive event perception [3] and extending the theory of image compositionality [18, 23, 41], we decompose temporal reasoning into atomic *Events*, parametric *Event Attributes*, and *Structural Logic*.

3. The TimeBlind Benchmark

In this section, we describe the design of TimeBlind, a diagnostic benchmark for rigorously evaluating compositional spatio-temporal understanding. Our design is grounded in two principles: (i) a temporal minimal-pair protocol that minimizes visual shortcuts and language priors, and (ii) a cognition-inspired compositionality taxonomy that decomposes temporal reasoning into three primitives: *Events*, *Event Attributes*, and *Structural Event Logic*.

3.1. Task Formulation

We structure each instance as $\mathcal{T} = (v_1, v_2, q_1, q_2)$, following the design below. This design creates a discriminative challenge where shortcut solutions are effectively “canceled out,” forcing the model to rely on temporal evidence.

Temporal Minimal Pairs. The video pair (v_1, v_2) constitutes a minimal pair: videos share identical static content (e.g., objects, background) but differ along an isolated temporal axis in our taxonomy (Section 3.2). Consequently, a model cannot distinguish v_1 from v_2 without explicitly modeling the temporal features (e.g., *opening* vs. *closing*).

Logical Complementarity. The questions (q_1, q_2) are complementary: for any question, the ground-truth answer flips between the two videos (e.g., $Ans(v_1, q_j) \neq Ans(v_2, q_j)$ for $j \in \{1, 2\}$). This design neutralizes language priors,

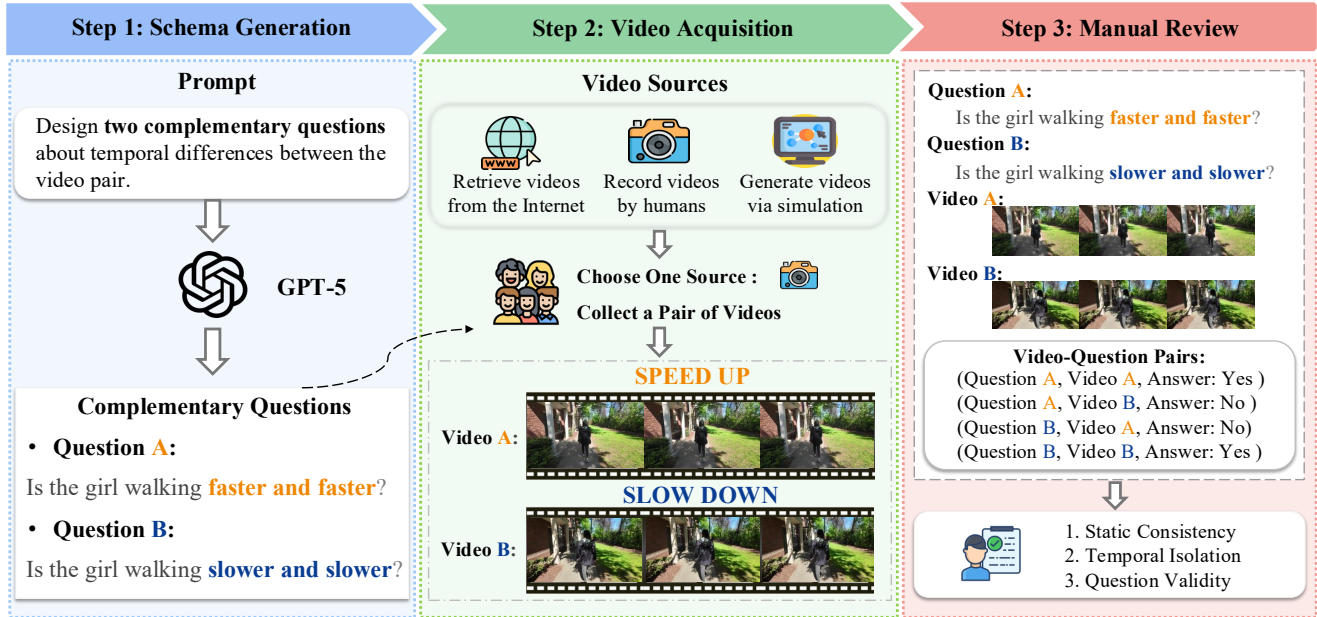


Figure 3. **Overview of the TimeBlind data construction pipeline.** **Stage 1 (Schema Generation):** We prompt GPT-5 to generate paired complementary questions targeting temporal differences. **Stage 2 (Video Acquisition):** We collect one video pair that matches the generated schema from one of the following sources: (i) Retrieving videos from the internet, (ii) Recording videos with humans, or (iii) Generating videos via simulation (e.g., Unity). We then pair these videos with the questions to form a candidate TimeBlind instance. **Stage 3 (Manual Review):** Human annotators manually review each instance to ensure: (i) *Static Consistency* (videos share identical static content), (ii) *Temporal Minimality* (the pair differs only in the targeted temporal factor), and (iii) *Question Validity* (QA pairs are clear and correct).

preventing models from exploiting textual plausibility.

Diagnostic Metrics. Following rigorous diagnostic protocols [23, 41], we report a hierarchy of metrics to mitigate shortcuts and assess the models’ genuine temporal understanding:

- **Accuracy (Acc):** Standard accuracy computed over all individual video–question trials.
- **Video Accuracy (V-Acc):** Measures visual consistency; correct only if the model answers *both* questions correctly for a single video v_j .
- **Question Accuracy (Q-Acc):** Measures textual consistency; correct only if the model answers question q_i correctly for *both* videos.
- **Instance Accuracy (I-Acc):** An instance is correct *if and only if* the model solves all four trials. I-Acc is our primary proxy for spatio-temporal understanding, as it necessitates reliably distinguishing the temporal difference between the paired videos.

3.2. A Hierarchy of Temporal Composition

To systematically evaluate temporal understanding, we organize TimeBlind around a compositional taxonomy mirroring a bottom-up cognitive process: from recognizing atomic changes (*Events*), to measuring their intrinsic properties (*Event Attributes*), and finally reasoning about their interdependencies (*Structural Event Logic*), shown in Figure 2.

I. Events. This foundational primitive tests the detection of atomic visual changes, demanding that models distinguish actual temporal evolution from static object presence.

- **Fine-Grained Actions:** Discriminating actions sharing visual context but differing in temporal dynamics (e.g., *opening* vs. *closing*) or semantic nuance (e.g., *placing* vs. *pushing*).
- **State Transitions:** Tracking object property changes over time (shape, size, color, emotion)—e.g., liquid changing from clear to opaque.

II. Event Attributes. This primitive tests the perception of *how* an event unfolds, requiring sensitivity to continuous or qualitative parameters rather than categorical recognition.

- **Kinematics:** Properties governing motion in space-time, including *Speed* (fast/slow), *Direction* (forward/backward), *Duration* (brief/long), and *Repetition* (once/twice).
- **Dynamics & Manner:** Properties reflecting physical forces or execution style, including *Force* (gentle/forceful) and *Magnitude* (slight/broad movement).

III. Structural Event Logic. This level tests *how* multiple events relate and compose into higher-order structures.

- **Temporal Topology:** We adopt Allen’s Interval Algebra [1] to comprehensively evaluate interval relations. Unlike prior work limited to simple sequencing (*before/after*), we cover all 13 relations in our videos, such as *overlaps*, *meets*, *starts*, and *during*.

- **Causal Contingency:** Identifying causal links between events (e.g., distinguishing *causation* from *correlation*).
- **Cross-Event Comparison:** Comparative reasoning across distinct events (e.g., “Does the person hold the cup *longer* than the book?”).

3.3. Data Construction Pipeline

We use a three-stage pipeline to collect high-quality TimeBlind data, as shown in Figure 3.

Step 1: Schema Generation. We prompt frontier LLMs (e.g., GPT-5) with the taxonomy definitions in Section 3.2 to generate structured specifications. Each output includes a pair of video descriptions (e.g., “*a girl walks in a park and gradually speeds up*” and “*a girl walks in a park and gradually slows down*”) that differ only in the target temporal factor (e.g., “*Speed*”), along with complementary questions (e.g., “*Is the girl walking faster and faster?*” and “*Is the girl walking slower and slower?*”).

Step 2: Video Acquisition. We collect video pairs from three sources: (i) internet retrieval (extracting temporally distinct segments from the same source video), (ii) human recording, or (iii) simulation (e.g., Unity) for precise temporal control. We then pair the collected videos with the questions to form a candidate TimeBlind instance.

Step 3: Rigorous Human Verification. Every instance undergoes strict manual review. Annotators verify (i) *Static Consistency*, ensuring both videos share identical static content; (ii) *Temporal Minimality*, confirming the pair differs only in the targeted temporal factor; and (iii) *Question Validity*, ensuring the QA pairs are clear and correct.

Statistics. TimeBlind comprises 2,400 curated video-question pairs, consisting of binary and two-choice multiple-choice questions. The video sources are distributed as follows: 24.0% from Internet Retrieval, 57.7% from Human Recording, and 18.3% from Simulation. For category distribution, we maintain an approximate 1:1:1 ratio across *Event*, *Event Attribute*, and *Structural Event Logic*. Within the *Temporal Topology* sub-category, we ensure a balanced distribution across 13 distinct Allen interval types.

4. Experimental Setup

Baseline Models. To thoroughly assess the challenges posed by TimeBlind, we evaluate over 20 MLLMs, including both frontier proprietary models (e.g., GPT-5 [32], Gemini 3 Pro [35], Claude [2], and Qwen3-VL Plus [4]) and state-of-the-art open-source models (e.g., Qwen3-VL [4], InternVL 3.5 [42], Molmo2 [9], PLM [8], Keye-VL [40], MiniCPM-V-4.5 [52], GLM-4.1V [12], and Eagle 2.5 [6]).

Evaluation Metrics. We adopt the hierarchy of metrics defined in Section 3.1, including Standard Accuracy (Acc), Video Accuracy (V-Acc), Question Accuracy (Q-Acc), and our primary metric, Instance Accuracy (I-Acc).

Implementation Details. We conduct most experiments on 4 NVIDIA H100 GPUs, each with 96GB of memory. For proprietary models (e.g., GPT-5, Gemini 3 Pro) and large-scale open-source models (e.g., Qwen3-VL-235B), we utilize the official APIs. Unless otherwise specified, we sample videos uniformly at 1 FPS (the default frame rate for most models) and evaluate models in a zero-shot setting. We also report experimental results across FPS settings from 1 to 10 in the Appendix 7.5.

5. Experimental Results

5.1. Main Results

In Table 2, we present the overall model performances on our benchmark, revealing that all models perform poorly in fine-grained temporal video understanding. While many models achieve relatively high standard accuracy (Acc), they struggle to achieve good Instance Accuracy (I-Acc), which requires correctly distinguishing both videos in a pair. For instance, despite leading models like GPT-5 and Gemini 3 Pro achieving high Acc scores of 77.3% and 76.2% respectively, they only reach 46.3% and 48.2% on I-Acc. This discrepancy indicates that models remain weak in temporal understanding and suggests that high performance on Acc is often driven by shortcuts, rather than a true understanding of temporal dynamics. Another interesting finding is that Question Accuracy (Q-Acc) is consistently lower than Video Accuracy (V-Acc) across all evaluated models. This indicates that models are more prone to hallucinating answers based on textual patterns than misinterpreting visual cues. From the results, we also observe that among open-source models, Molmo2-8B—despite its smaller size—outperforms a range of sub-10B models and even surpasses the much larger Qwen3-VL-235B by 5.4% on I-Acc, becoming the leading open-source model. This suggests that with effective design choices in model architecture, data curation, and training process, even smaller models can capture fine-grained temporal dynamics. Lastly, our results show a significant gap between open-source and proprietary models in understanding fine-grained temporal dynamics. Specifically, Molmo2-8B still lags substantially behind GPT-5 and Gemini 3 Pro, with gaps of 15.1% and 17.0% I-Acc, respectively.

Human Evaluation. We validate TimeBlind with four independent annotators (Table 2). Each annotator saw only one question and one video at a time, with the four (v, q) pairs from each instance distributed across different annotators. Humans achieved 98.2% I-Acc—exceeding Gemini 3 Pro by 50%—showing that temporal dynamics in TimeBlind are clear to humans while remaining challenging for MLLMs.

5.2. Category-Wise Diagnosis

To pinpoint specific cognitive deficits, we analyze performance across the TimeBlind taxonomy in Table 3.

Table 2. **Main Results on TimeBlind.** We use uniform sampling at 1 FPS and evaluate all models using default configurations. Metrics are reported following Section 3, with I-Acc as the primary metric. The table is divided into *Open-Source* (grouped by size: $< 10B$ and $> 10B$), and *Closed-Source models*. Our results suggest that all models perform poorly on fine-grained temporal video understanding, with none of the methods achieving over 50% I-Acc. Best results are shown in **bold**.

Model	Size	Q-Acc (%)	V-Acc (%)	Acc (%)	I-Acc (%)
Random Chance	-	25.0	25.0	50.0	6.3
Human Evaluation	-	98.8	98.9	99.3	98.2
Open-Source Models					
<i>Model Size $\leq 10B$</i>					
Keye-VL-1.5 [40]	8B	17.5	26.5	55.6	6.6
LLaVA-Video [55]	7B	18.1	34.8	57.8	9.6
InternVL3.5 [42]	8B	30.9	39.0	59.3	13.3
PLM-8B [8]	8B	24.4	32.3	60.6	13.9
F-16 [26]	7B	23.5	39.7	60.5	14.5
video-SALMONN 2+ [39]	7B	26.3	41.7	61.7	16.2
Qwen3-VL [4]	4B	29.9	43.7	62.6	17.7
MiniCPM-V-4.5 [52]	8B	31.7	43.6	62.7	18.3
VideoChat-Flash [25]	7B	31.7	45.4	64.9	19.5
Qwen3-VL [4]	8B	34.5	43.6	64.8	19.7
GLM-4.1V [12]	9B	32.3	45.2	62.0	19.7
Eagle2.5 [6]	8B	32.2	45.7	64.7	20.2
Molmo2 [9]	8B	41.0	52.4	68.7	31.2
<i>Model Size $\geq 30B$</i>					
InternVL 3.5 [42]	30B	31.4	41.3	61.1	17.0
LLaVA-Video [55]	72B	28.8	45.0	63.2	18.5
Qwen3-VL [4]	235B	38.1	51.7	66.9	25.8
Closed-Source Models					
Claude Sonnet 4.5 [2]	unknown	30.6	38.5	59.2	13.6
Qwen3-VL Plus [4]	unknown	38.7	52.6	67.3	26.0
GPT-5 mini [32]	unknown	42.3	53.7	68.3	30.0
Gemini 2.5 Flash [17]	unknown	47.6	53.7	68.9	33.2
Gemini 2.5 Pro [17]	unknown	56.0	62.6	74.2	43.5
GPT-5 [32]	unknown	58.5	67.2	77.3	46.3
Gemini 3 Pro [35]	unknown	60.2	66.0	76.2	48.2

- Performance Across Hierarchical Levels.** We observe notable performance gaps across different temporal understanding tasks, with models generally performing better on discrete *Events* (e.g., distinguishing atomic actions) compared to continuous *Event Attributes* or *Structural Event Logic*. For example, GPT-5 achieves 58.3% accuracy in the *Event* category (peaking at 62.5% for *Fine-Grained Action*), but performance declines to 32.3% for *Event Attributes* and 48.4% for *Structural Event Logic*. Furthermore, we note that while LLaVA-Video-72B and InternVL 3.5-38B achieve moderate performance on *Events* (32.4% and 33.3%), they perform only slightly above random chance in the other two categories.

- Physical Dynamics.** Performance is lowest in the *Event Attributes* category, which requires sensitivity to continuous or qualitative parameters such as kinematics (speed) and dynamics (force, magnitude). Strong proprietary models GPT-5 and Gemini 3 Pro only achieve 32.3% and 36.7% I-Acc, respectively, while the top-performing open-source model, Molmo2-8B, scores only 20.3%. Furthermore, several models, including Qwen3-VL-235B, LLaVA-Video-72B, and InternVL 3.5-38B, perform at near random chance, struggling to distinguish nuances such as *gentle vs. forceful*, *fast vs. slow*, or *large vs. small amplitude*. These results expose a systematic deficiency in current models' understanding of low-level, physics-

Table 3. **Category-Wise I-Acc (%) on TimeBlind.** This table reports I-Acc for advanced models across 11 fine-grained temporal understanding tasks to pinpoint specific cognitive deficits. Due to space constraints, we use the following abbreviations for temporal categories: *FG Action*: Fine-Grained Action, *State Trans*: State Transitions, *Mag*: Magnitude, *Dir*: Direction, *Dur*: Duration, *Rep*: Repetition, *Temp Topo*: Temporal Topology, *Causal Cont*: Causal Contingency, and *Cross Comp*: Cross-Event Comparison. We include the overall performance (*Avg*) for each high-level category alongside the fine-grained categories. The results show clear performance gaps across categories. Models generally perform well on discrete *Events*, but struggle with *Event Attributes*, which require low-level physical understanding such as *Speed* and *Force*. The best results are **bolded**.

Models	Event			Event Attribute						Structural Event Logic				
	FG Action	State Trans	Avg	Speed	Force	Mag	Dir	Dur	Rep	Avg.	Temp Topo	Causal Cont	Cross Comp	Avg
Random Chance	6.3	6.3	6.3	6.3	6.3	6.3	6.3	6.3	6.3	6.3	6.3	6.3	6.3	6.3
Open-Source Models														
InternVL3.5-38B	30.7	33.6	32.4	3.6	0.0	4.1	15.7	50.0	5.6	9.4	10.2	5.0	11.4	9.4
LLaVA-Video-72B	23.9	40.5	33.3	0.0	5.6	0.0	8.6	40.0	11.1	7.3	10.2	10.0	18.2	12.0
Qwen3-VL-235B-Instruct	38.6	47.4	43.6	3.6	5.6	8.3	18.6	75.0	13.9	12.5	20.4	7.5	25.0	18.8
Molmo2-8B	54.5	49.1	52.0	10.7	11.1	12.5	31.4	25.0	19.4	20.3	18.5	15.0	18.2	17.7
Closed-Source Models														
GPT-5	62.5	55.2	58.3	21.4	16.7	25.0	40.0	75.0	30.6	32.3	49.1	50.0	45.5	48.4
Gemini 3 Pro	49.4	49.0	49.2	58.3	50.0	45.0	23.1	75.0	26.5	36.7	69.0	37.5	50.0	57.8

related temporal dynamics.

- **Gap in Structural Logic.** We observe that GPT-5 and Gemini 3 Pro demonstrate relatively robust performance in *Structural Logic Analysis*, achieving 48.4% and 57.8%, respectively. However, a striking disparity emerges between these proprietary models and their open-source counterparts. The leading Qwen3-VL-235B achieves only 18.8%—trailing Gemini 3 Pro by 39.0%—and scores a mere 7.5% in the *Causal Contingency* sub-category. These results indicate that while open-source MLLMs may effectively detect isolated events, they lack the necessary logical abilities to reason about precise relationships between multiple events, particularly regarding causality.

5.3. Shortcut Analysis

To verify that solving TimeBlind requires genuine temporal understanding, we report three tests with GPT-5 in Table 4.

- **Single-Frame Bias.** This experiment evaluates whether TimeBlind requires reasoning over a sequence rather than exploiting information from a single static frame. In this setting, the model is provided with the question and only one randomly sampled frame. The results show that GPT-5 performs poorly (4.5% I-Acc), indicating that TimeBlind requires sequential modeling.
- **Language-Only Bias.** This experiment evaluates the influence of language priors within TimeBlind. In this setting, the model is provided with only the question without any visual information. As shown in Table 4, GPT-5 achieves only 1.5% I-Acc, demonstrating that visual information is essential for our benchmark setting.
- **Visual-Cue Shortcuts.** Finally, we examine whether the

Table 4. **Shortcut Analysis.** We report I-Acc for three shortcut baselines using GPT-5 with 1 FPS sampling: Single Frame Bias (a question paired with a randomly selected video frame), Language Only (a question without visual input), and Visual-Cue (a question with shuffled video frames). The results demonstrate that solving TimeBlind requires genuine temporal understanding.

Setting	Acc	Q-Acc	V-Acc	I-Acc
Random Chance	50.0	25.0	25.0	6.3
GPT-5 (<i>baseline</i>)	77.3	58.5	67.2	46.3
Single Frame	52.2	14.8	32.3	4.5
Language Only	47.3	9.4	26.6	1.5
Visual-Cue	49.6	10.7	23.2	3.0

benchmark can be solved by exploiting static visual cues that happen to correlate with temporal dynamics. In this setting, we shuffle the order of the input frames sampled at 1 FPS to ensure that successfully completing the task requires an understanding of the temporal sequence rather than merely detecting specific objects. GPT-5 achieves only 3.0% I-Acc in this setting, indicating that solving TimeBlind requires strict temporal understanding.

Moreover, Acc scores across all three settings hover around random chance, suggesting that correctly answering even a single question in TimeBlind requires strict temporal understanding. In summary, TimeBlind stands as a temporal-centric benchmark robust against shortcut solutions.

Table 5. **Effect of Model Size and Input Frames.** We report the I-Acc performance of several open-source and proprietary models across various model sizes and input frame counts. The results indicate that simply scaling up model size and frame count fails to significantly improve performance. Best results within each model grouping are highlighted in **bold** with a gray background.

Size	Frames	Q-Acc	V-Acc	Acc	I-Acc
<i>InternVL 3.5</i>					
8B	8	34.7	41.5	61.9	17.1
	16	35.6	42.7	62.6	18.5
	32	35.6	43.1	62.9	17.4
14B	8	34.2	46.3	63.6	19.9
	16	38.2	47.6	65.2	23.5
	32	37.7	49.4	65.4	23.9
38B	8	35.8	42.8	63.2	20.4
	16	38.9	45.4	64.5	24.4
	32	39.4	46.4	65.9	25.1
<i>LLaVA-Video</i>					
7B	8	21.8	37.9	59.5	12.4
	16	23.4	39.8	60.1	13.9
	32	24.7	41.1	60.7	14.6
72B	8	29.7	46.2	63.5	19.9
	16	32.4	48.4	65.0	23.0
	32	32.8	48.7	65.2	23.3
<i>GPT-5</i>					
unknown	8	60.9	69.2	79.0	49.1
	16	62.7	69.9	79.5	50.9
	32	61.1	68.2	78.7	48.3

Table 6. **Inference-Time Reasoning.** We compare standard Qwen3-VL and GPT-5 models against their reasoning-enabled *Thinking* counterparts across various settings, using the I-Acc metric at 1 FPS sampling. Results show that sufficient inference-time reasoning depth improves temporal understanding, but is far from sufficient to solve TimeBlind.

Model	Mode	I-Acc (%)	Δ
<i>Qwen3-VL-8B</i>	Standard	19.6	-
	Thinking	27.8	+8.2
<i>Qwen3-VL-235B</i>	Standard	25.8	-
	Thinking	36.3	+10.4
<i>GPT-5</i>	Standard	46.3	-
	Low-Thinking	43.8	-2.5
	Med-Thinking	47.9	+1.7
	High-Thinking	49.6	+3.3

5.4. Additional Analysis

We conduct ablation studies on three factors: input frames, model size, and inference-time reasoning. We also report experimental results across FPS settings from 1 to 10 in the

Appendix 7.5.

Number of Input Frames. Table 5 shows that increasing frames from 8 to 32 yields only marginal I-Acc gains: 1–5% for InternVL 3.5 and LLaVA-Video, and less than 2% for GPT-5. Even with sufficient frames, models struggle with TimeBlind, revealing fundamental limitations in fine-grained temporal understanding.

Impact of Model Size. Despite an $10\times$ parameter increase for LLaVA-Video (7B to 72B) and $5\times$ for InternVL 3.5 (8B to 38B), both show less than 10% I-Acc improvement across all frame settings (Table 5). Simply scaling model size does not yield robust spatio-temporal understanding.

Test-Time Scaling (Reasoning). Table 6 shows inference-time reasoning results. Due to computational cost, we evaluate on a random 30% subset of the data. For Qwen3-VL, *Thinking* variants outperform *Instruct* models, with the 235B model gaining 10.4% but still achieving only 36.3% I-Acc. For GPT-5, varying reasoning effort from Low to High yields modest gains, peaking at 49.6% I-Acc—still far below human performance (98.2%). While deeper reasoning helps, it remains insufficient to solve TimeBlind.

6. Conclusion

We introduce TimeBlind, a diagnostic benchmark designed to rigorously assess the compositional spatio-temporal reasoning capabilities of MLLMs. By organizing evaluation around a structured cognitive taxonomy—spanning atomic *events*, *event attribute*, and *structural event logic*—and employing a strict minimal-pair design, we effectively isolate temporal understanding from static and linguistic shortcuts. Our evaluation reveals a significant gap: despite rapid progress in static vision-language tasks, current state-of-the-art models remain largely “time-blind,” trailing human performance by 50% and exhibiting fragility in fine-grained event attributes and temporal logic. We believe that TimeBlind will be beneficial for developing MLLMs capable of understanding fine-grained temporal dynamics.

Impact Statement

This work introduces a diagnostic benchmark for evaluating temporal understanding in video-language models. We anticipate several positive impacts: (1) TimeBlind can guide development of more temporally-aware models, which is critical for applications in robotics, autonomous driving, and assistive technologies; (2) by revealing systematic failure modes, our benchmark may help practitioners avoid deploying models in safety-critical scenarios where temporal reasoning is essential.

We also acknowledge potential concerns. Improved video understanding capabilities could enable more sophisticated surveillance or content moderation systems, raising privacy considerations. Additionally, our benchmark primarily fea-

tures videos from controlled settings and internet sources, which may not represent the full diversity of real-world scenarios or populations. We encourage future work to expand temporal reasoning evaluation to more diverse contexts. The benchmark itself does not enable direct harmful applications, as it evaluates rather than enhances model capabilities. We will release our data and code to support reproducible research.

Acknowledgments

This work was supported by the National Institutes of Health Award 1R01HD111074-01. We thank Zhiqiu Lin, Kewen Wu, and Deva Ramanan for their invaluable discussions during the development of this work.

References

- [1] James F Allen. Maintaining knowledge about temporal intervals. *Communications of the ACM*, 26(11):832–843, 1983. 2, 4
- [2] Anthropic. Introducing claude sonnet 4.5, 2025. Anthropic News. 5, 6
- [3] Emmon Bach. The algebra of events. *Linguistics and philosophy*, pages 5–16, 1986. 3
- [4] Shuai Bai, Yuxuan Cai, Ruizhe Chen, Keqin Chen, Xionghui Chen, Zesen Cheng, Lianghao Deng, Wei Ding, Chang Gao, Chunjiang Ge, Wenbin Ge, Zhifang Guo, Qidong Huang, Jie Huang, Fei Huang, Binyuan Hui, Shutong Jiang, Zhaohai Li, Mingsheng Li, Mei Li, Kaixin Li, Zicheng Lin, Junyang Lin, Xuejing Liu, Jiawei Liu, Chenglong Liu, Yang Liu, Dayiheng Liu, Shixuan Liu, Dunjie Lu, Ruilin Luo, Chenxu Lv, Rui Men, Lingchen Meng, Xuancheng Ren, Xingzhang Ren, Sibao Song, Yuchong Sun, Jun Tang, Jianhong Tu, Jianqiang Wan, Peng Wang, Pengfei Wang, Qiuyue Wang, Yuxuan Wang, Tianbao Xie, Yiheng Xu, Haiyang Xu, Jin Xu, Zhibo Yang, Mingkun Yang, Jianxin Yang, An Yang, Bowen Yu, Fei Zhang, Hang Zhang, Xi Zhang, Bo Zheng, Humen Zhong, Jingren Zhou, Fan Zhou, Jing Zhou, Yuanzhi Zhu, and Ke Zhu. Qwen3-vl technical report, 2025. 1, 2, 5, 6
- [5] Mu Cai, Reuben Tan, Jianrui Zhang, Bocheng Zou, Kai Zhang, Feng Yao, Fangrui Zhu, Jing Gu, Yiwu Zhong, Yuzhang Shang, et al. Temporalbench: Benchmarking fine-grained temporal understanding for multimodal video models. *arXiv preprint arXiv:2410.10818*, 2024. 3
- [6] Guo Chen, Zhiqi Li, Shihao Wang, Jindong Jiang, Yicheng Liu, Lidong Lu, De-An Huang, Wonmin Byeon, Matthieu Le, Max Ehrlich, Tong Lu, Limin Wang, Bryan Catanzaro, Jan Kautz, Andrew Tao, Zhiding Yu, and Guilin Liu. Eagle 2.5: Boosting long-context post-training for frontier vision-language models. In *NeurIPS*, 2025. 5, 6
- [7] Junhao Cheng, Yuying Ge, Teng Wang, Yixiao Ge, Jing Liao, and Ying Shan. Video-holmes: Can mllm think like holmes for complex video reasoning? *arXiv preprint arXiv:2505.21374*, 2025. 2
- [8] Jang Hyun Cho, Andrea Madotto, Effrosyni Mavroudi, Triantafyllos Afouras, Tushar Nagarajan, Muhammad Maaz, Yale Song, Tengyu Ma, Shuming Hu, Suyog Jain, et al. Perceptionlm: Open-access data and models for detailed visual understanding. *arXiv preprint arXiv:2504.13180*, 2025. 1, 2, 5, 6
- [9] Christopher Clark, Jieyu Zhang, Zixian Ma, Jae Sung Park, Mohammadreza Salehi, Rohun Tripathi, Sangho Lee, Zhongzheng Ren, Chris Dongjoo Kim, Yinuo Yang, Vincent Shao, Yue Yang, Weikai Huang, Ziqi Gao, Taira Anderson, Jianrui Zhang, Jitesh Jain, George Stoica, Winson Han, Ali Farhadi, and Ranjay Krishna. Molmo2: Open weights and data for vision-language models with video understanding and grounding. Technical report, Allen Institute for AI, 2025. 1, 2, 5, 6
- [10] Chaoyou Fu, Yuhan Dai, Yongdong Luo, Lei Li, Shuhuai Ren, Renrui Zhang, Zihan Wang, Chenyu Zhou, Yunhang Shen, Mengdan Zhang, et al. Video-mme: The first-ever comprehensive evaluation benchmark of multi-modal llms in video analysis. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 24108–24118, 2025. 1, 2
- [11] Xingyu Fu, Yushi Hu, Bangzheng Li, Yu Feng, Haoyu Wang, Xudong Lin, Dan Roth, Noah A Smith, Wei-Chiu Ma, and Ranjay Krishna. Blink: Multimodal large language models can see but not perceive. In *European Conference on Computer Vision*, pages 148–166. Springer, 2024. 3
- [12] Wenyi Hong, Wenmeng Yu, Xiaotao Gu, Guo Wang, Guobing Gan, Haomiao Tang, Jiale Cheng, Ji Qi, Junhui Ji, Lihang Pan, et al. Glm-4.1 v-thinking: Towards versatile multimodal reasoning with scalable reinforcement learning. *arXiv preprint arXiv:2507.01006*, 2025. 5, 6
- [13] Kairui Hu, Penghao Wu, Fanyi Pu, Wang Xiao, Yuanhan Zhang, Xiang Yue, Bo Li, and Ziwei Liu. Video-mmmu: Evaluating knowledge acquisition from multi-discipline professional videos. *arXiv preprint arXiv:2501.13826*, 2025. 1, 2
- [14] Yunseok Jang, Yale Song, Youngjae Yu, Youngjin Kim, and Gunhee Kim. Tgif-qa: Toward spatio-temporal reasoning in visual question answering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2758–2766, 2017. 1
- [15] Yunseok Jang, Yale Song, Youngjae Yu, Youngjin Kim, and Gunhee Kim. Tgif-qa: Toward spatio-temporal reasoning in visual question answering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2758–2766, 2017. 2
- [16] Justin Johnson, Bharath Hariharan, Laurens Van Der Maaten, Li Fei-Fei, C Lawrence Zitnick, and Ross Girshick. Clevr: A diagnostic dataset for compositional language and elementary visual reasoning. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2901–2910, 2017. 3
- [17] Koray Kavukcuoglu. Gemini 2.5: Our most intelligent ai model, 2025. Google Keyword Blog. Last updated March 26, 2025. 6
- [18] Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A Shamma, et al. Visual genome: Connecting language and vision using crowdsourced dense image annotations. *International journal of computer vision*, 123(1):32–73, 2017. 2, 3
- [19] Benno Krojer, Mojtaba Komeili, Candace Ross, Quentin Garrido, Koustuv Sinha, Nicolas Ballas, and Mahmoud Assran. A shortcut-aware video-qa benchmark for physical understanding via minimal video pairs. *arXiv preprint arXiv:2506.09987*, 2025. 1, 2, 3
- [20] Jie Lei, Licheng Yu, Mohit Bansal, and Tamara Berg. Tvqa: Localized, compositional video question answering. In *Proceedings of the 2018 conference on empirical methods in natural language processing*, pages 1369–1379, 2018. 2
- [21] Jie Lei, Tamara Berg, and Mohit Bansal. Revealing single frame bias for video-and-language learning. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 487–507, 2023. 1, 2

- [22] Baiqi Li, Zhiqiu Lin, Deepak Pathak, Jiayao Li, Yixin Fei, Kewen Wu, Tiffany Ling, Xide Xia, Pengchuan Zhang, Graham Neubig, et al. Genai-bench: Evaluating and improving compositional text-to-visual generation. *arXiv preprint arXiv:2406.13743*, 2024. 2
- [23] Baiqi Li, Zhiqiu Lin, Wenxuan Peng, Jean de Dieu Nyandwi, Daniel Jiang, Zixian Ma, Simran Khanuja, Ranjay Krishna, Graham Neubig, and Deva Ramanan. Naturalbench: Evaluating vision-language models on natural adversarial samples. *Advances in Neural Information Processing Systems*, 37:17044–17068, 2024. 1, 2, 3, 4
- [24] Kunchang Li, Yali Wang, Yinan He, Yizhuo Li, Yi Wang, Yi Liu, Zun Wang, Jilan Xu, Guo Chen, Ping Luo, et al. Mvbench: A comprehensive multi-modal video understanding benchmark. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22195–22206, 2024. 1, 2
- [25] Xinhao Li, Yi Wang, Jiashuo Yu, Xiangyu Zeng, Yuhan Zhu, Haiyan Huang, Jianfei Gao, Kunchang Li, Yinan He, Chenting Wang, et al. Videochat-flash: Hierarchical compression for long-context video modeling. *arXiv preprint arXiv:2501.00574*, 2024. 6
- [26] Yixuan Li, Changli Tang, Jimin Zhuang, Yudong Yang, Guangzhi Sun, Wei Li, Zejun Ma, and Chao Zhang. Improving llm video understanding with 16 frames per second. *arXiv preprint arXiv:2503.13956*, 2025. 6
- [27] Yuanxin Liu, Shicheng Li, Yi Liu, Yuxiang Wang, Shuhuai Ren, Lei Li, Sishuo Chen, Xu Sun, and Lu Hou. Tempcompass: Do video llms really understand videos? *arXiv preprint arXiv:2403.00476*, 2024. 1, 2, 3
- [28] Xinwei Long, Kai Tian, Peng Xu, Guoli Jia, Jingxuan Li, Sa Yang, Yihua Shao, Kaiyan Zhang, Che Jiang, Hao Xu, Yang Liu, Jiaheng Ma, and Bowen Zhou. Adsqa: Towards advertisement video understanding. In *ICCV*, 2025. 2
- [29] Wentao Ma, Weiming Ren, Yiming Jia, Zhuofeng Li, Ping Nie, Ge Zhang, and Wenhui Chen. Videoeval-pro: Robust and realistic long video understanding evaluation. *arXiv preprint arXiv:2505.14640*, 2025. 2
- [30] Kartikeya Mangalam, Raiymbek Akshulakov, and Jitendra Malik. Egoschema: A diagnostic benchmark for very long-form video language understanding. *Advances in Neural Information Processing Systems*, 36:46212–46244, 2023. 2
- [31] Arsha Nagrani, Sachit Menon, Ahmet Iscen, Shyamal Buch, Ramin Mehran, Nilpa Jha, Anja Hauth, Yukun Zhu, Carl Vondrick, Mikhail Sirotenko, et al. Minerva: Evaluating complex video reasoning. *arXiv preprint arXiv:2505.00681*, 2025. 2
- [32] OpenAI. GPT-5 system card, 2025. Accessed: 2026-01-05. 1, 2, 5, 6
- [33] Yulu Pan, Ce Zhang, and Gedas Bertasius. Basket: A large-scale video dataset for fine-grained skill estimation. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 28952–28962, 2025. 1, 2
- [34] Viorica Patraucean, Lucas Smaira, Ankush Gupta, Adria Recasens, Larisa Markeeva, Dylan Banarse, Skanda Koppula, Mateusz Malinowski, Yi Yang, Carl Doersch, et al. Perception test: A diagnostic benchmark for multimodal video models. *Advances in Neural Information Processing Systems*, 36:42748–42761, 2023. 2
- [35] Sundar Pichai, Demis Hassabis, and Koray Kavukcuoglu. A new era of intelligence with gemini 3. The Keyword (Google Blog), 2025. Accessed: 2026-01-05. 1, 2, 5, 6
- [36] Ziyao Shangguan, Chuhan Li, Yuxuan Ding, Yanan Zheng, Yilun Zhao, Tesca Fitzgerald, and Arman Cohan. Tomato: Assessing visual temporal reasoning capabilities in multimodal foundation models. *arXiv preprint arXiv:2410.23266*, 2024. 2, 3
- [37] Enxin Song, Wenhao Chai, Guan hong Wang, Yucheng Zhang, Haoyang Zhou, Feiyang Wu, Haozhe Chi, Xun Guo, Tian Ye, Yanting Zhang, et al. Moviechat: From dense token to sparse memory for long video understanding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18221–18232, 2024. 2
- [38] Enxin Song, Wenhao Chai, Weili Xu, Jianwen Xie, Yuxuan Liu, and Gaoang Wang. Video-mmlu: A massive multi-discipline lecture understanding benchmark. *arXiv preprint arXiv:2504.14693*, 2025. 2
- [39] Yunlong Tang, Jing Bi, Siting Xu, Luchuan Song, Susan Liang, Teng Wang, Daoan Zhang, Jie An, Jingyang Lin, Rongyi Zhu, et al. Video understanding with large language models: A survey. *IEEE Transactions on Circuits and Systems for Video Technology*, 2025. 1, 6
- [40] Kwai Keye Team, Biao Yang, Bin Wen, Changyi Liu, Chenglong Chu, Chengru Song, Chongling Rao, Chuan Yi, Da Li, Dunju Zang, et al. Kwai keye-v1 technical report. *arXiv preprint arXiv:2507.01949*, 2025. 5, 6
- [41] Tristan Thrush, Ryan Jiang, Max Bartolo, Amanpreet Singh, Adina Williams, Douwe Kiela, and Candace Ross. Winoground: Probing vision and language models for visiolinguistic compositionality. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5238–5248, 2022. 2, 3, 4
- [42] Weiyun Wang, Zhangwei Gao, Lixin Gu, Hengjun Pu, Long Cui, Xingguang Wei, Zhaoyang Liu, Linglin Jing, Shenglong Ye, Jie Shao, et al. Internv13. 5: Advancing open-source multimodal models in versatility, reasoning, and efficiency. *arXiv preprint arXiv:2508.18265*, 2025. 5, 6
- [43] Wei han Wang, Zehai He, Wenyi Hong, Yean Cheng, Xiaohan Zhang, Ji Qi, Ming Ding, Xiaotao Gu, Shiyu Huang, Bin Xu, et al. Lvbench: An extreme long video understanding benchmark. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 22958–22967, 2025. 1, 2
- [44] Bo Wu, Shoubin Yu, Zhenfang Chen, Joshua B Tenenbaum, and Chuang Gan. Star: A benchmark for situated reasoning in real-world videos. *arXiv preprint arXiv:2405.09711*, 2024. 2
- [45] Haoning Wu, Dongxu Li, Bei Chen, and Junnan Li. Longvideobench: A benchmark for long-context interleaved video-language understanding. *Advances in Neural Information Processing Systems*, 37:28828–28857, 2024. 2
- [46] Junbin Xiao, Xindi Shang, Angela Yao, and Tat-Seng Chua. Next-qa: Next phase of question-answering to explaining temporal actions. In *Proceedings of the IEEE/CVF conference on*

computer vision and pattern recognition, pages 9777–9786, 2021. [2](#)

- [47] Jun Xu, Tao Mei, Ting Yao, and Yong Rui. Msr-vtt: A large video description dataset for bridging video and language. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5288–5296, 2016. [2](#)
- [48] Zihui Xue, Mi Luo, and Kristen Grauman. Seeing the arrow of time in large multimodal models. *arXiv preprint arXiv:2506.03340*, 2025. [3](#)
- [49] Jingkang Yang, Shuai Liu, Hongming Guo, Yuhao Dong, Xiamengwei Zhang, Sicheng Zhang, Pengyun Wang, Zitang Zhou, Binzhu Xie, Ziyue Wang, et al. Egolife: Towards egocentric life assistant. *arXiv preprint arXiv:2503.03803*, 2025. [2](#)
- [50] Shusheng Yang, Jihan Yang, Pinzhi Huang, Ellis Brown, Zihao Yang, Yue Yu, Shengbang Tong, Zihan Zheng, Yifan Xu, Muhan Wang, et al. Cambrian-s: Towards spatial supersensing in video. *arXiv preprint arXiv:2511.04670*, 2025. [1](#)
- [51] Han Yi, Yulu Pan, Feihong He, Xinyu Liu, Benjamin Zhang, Oluwatuminu Oguntola, and Gedas Bertasius. Exact: A video-language benchmark for expert action analysis. *arXiv preprint arXiv:2506.06277*, 2025. [2](#)
- [52] Tianyu Yu, Zefan Wang, Chongyi Wang, Fuwei Huang, Wenshuo Ma, Zhihui He, Tianchi Cai, Weize Chen, Yuxiang Huang, Yuanqian Zhao, et al. Minicpm-v 4.5: Cooking efficient mllms via architecture, data, and training recipe. *arXiv preprint arXiv:2509.18154*, 2025. [5](#), [6](#)
- [53] Zhou Yu, Dejing Xu, Jun Yu, Ting Yu, Zhou Zhao, Yueting Zhuang, and Dacheng Tao. Activitynet-qa: A dataset for understanding complex web videos via question answering. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 9127–9134, 2019. [1](#), [2](#)
- [54] Jianrui Zhang, Mu Cai, and Yong Jae Lee. Vinoground: Scrutinizing lms over dense temporal reasoning with short videos. *arXiv preprint arXiv:2410.02763*, 2024. [2](#), [3](#)
- [55] Yuanhan Zhang, Jinming Wu, Wei Li, Bo Li, Zejun Ma, Ziwei Liu, and Chunyuan Li. Video instruction tuning with synthetic data. *arXiv preprint arXiv:2410.02713*, 2024. [6](#)
- [56] Junjie Zhou, Yan Shu, Bo Zhao, Boya Wu, Zhengyang Liang, Shitao Xiao, Minghao Qin, Xi Yang, Yongping Xiong, Bo Zhang, et al. Mlvu: Benchmarking multi-task long video understanding. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 13691–13701, 2025. [1](#), [2](#)
- [57] Yiyang Zhou, Linjie Li, Shi Qiu, Zhengyuan Yang, Yuyang Zhao, Siwei Han, Yangfan He, Kangqi Li, Haonian Ji, Zihao Zhao, et al. Glimpse: Do large vision-language models truly think with videos or just glimpse at them? In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 27830–27844, 2025. [3](#)

TimeBlind: A Spatio-Temporal Compositionality Benchmark for Video LLMs

Supplementary Material

7. Appendix

Outline

This document supplements the main paper by providing additional details on the TimeBlind benchmark, ablation studies, and representative failure cases. The contents are organized as follows:

- **The TimeBlind Benchmark:** data collection procedures, prompts used for data generation, quality control process, benchmark statistics, and taxonomy details.
- **Additional Ablation Studies:** additional ablation experiments on high-FPS settings and across video subsets with different levels of video similarity.
- **Impact Statement:** a discussion of the broader impact of the TimeBlind benchmark.
- **Failure Cases:** representative failure examples (1–7).

7.1. Data Collection Processes

Details on human-recorded videos. We recruited 10 participants and provided each with generated caption pairs and clear instructions: (1) record two videos matching the given captions; (2) maintain static consistency (same people, props, scene, camera angle); (3) ensure temporal minimality (videos differ only in the targeted temporal factor); and (4) meet quality standards (clear footage, similar duration within pairs). Videos are not limited to single-person scenes—some include multi-person interactions across settings such as parks, streets, restaurants, and offices.

Details on internet video collection. Annotators searched YouTube using LLM-generated captions or relevant keywords. They manually identified clips matching caption pairs. Where possible, both clips were sourced from the same source video to maximize static consistency.

Details on generated video collection. We used tools such as Unity to create videos for each caption pair. During generation, we ensured static visual consistency within each video pair, introduced differences only in the targeted temporal aspect, and kept the two videos in each pair nearly the same in duration.

7.2. Prompt Details.

In this section, we present the full prompt template used in our Human–AI collaborative pipeline (see Section 3.3). The system prompt for TimeBlind consists of three parts: (1) a **General Prompt** that defines the role and task, (2) **Detailed Requirements** that specify the core generation rules and logical constraints, and (3) an **Output Format** that

standardizes the response structure. The prompt includes strict structural constraints to ensure high-quality minimal pairs. When using this prompt, we provide a single category as the “Input” each time.

Part 1: Role, Context, and Task Definition

Role

You are an expert in dataset curation and benchmark design for Multimodal Large Language Models (MLLMs).

Context

We are constructing a benchmark called “**TimeBlind**” to evaluate the temporal perception capabilities of MLLMs. In TimeBlind, each data instance consists of two videos (v_1, v_2) and two questions (q_1, q_2).

- The videos are visually similar (same objects/background) but contain distinct temporal dynamics.
- The questions are designed to distinguish these specific temporal features.
- During testing, the model receives only **one video** and **one question** at a time.

Task

I will provide you with a specific “**Temporal Dynamic Category**”. Your task is to generate a valid data instance (one pair of videos and one pair of questions) that perfectly fits the requirements below.

Part 2: Generation Requirements and Logical Constraints

Requirements

1. Video Generation (v_1 and v_2)

- **Format:** Provide detailed text descriptions for both videos.
- **Feasibility:** The scenes must be easy to collect via web search, self-recording, or synthetic generation.
- **Visual Consistency:** v_1 and v_2 must share the same objects, background, and camera angle. They should look nearly identical in a static frame; the difference must be purely temporal.

2. Question Generation (q_1 and q_2)

- **Type:** Two-choice multiple-choice or Binary Questions (Yes/No).
- **Relevance:** Both questions must relate to the provided “Temporal Dynamic Category.”
- **Temporal Dependency:** The questions must

strictly require temporal understanding. They **cannot** be solvable by looking at a single static frame.

Logical Constraints (Crucial)

The answers (A) must adhere to the following logic matrix to ensure the model is truly perceiving time:

- **Constraint A (Cross-Video Difference):** The answer changes depending on the video.

$$A(q_1, v_1) \neq A(q_1, v_2) \quad \text{and} \quad A(q_2, v_1) \neq A(q_2, v_2)$$

- **Constraint B (Intra-Video Difference):** The two questions have different answers within the same video.

$$A(q_1, v_1) \neq A(q_2, v_1) \quad \text{and} \quad A(q_1, v_2) \neq A(q_2, v_2)$$

Part 3: Output Specifications

Output Format

Please present the result in the following structured format:

Temporal Category: [Input Category]

Video Descriptions:

- v_1 : [Description]
- v_2 : [Description]
- **Collection Source:** [e.g., Synthetic / Self-shot / Web]

Questions & Answers:

- q_1 : [Question Text]
Answer in v_1 : [Answer] / Answer in v_2 : [Answer]
- q_2 : [Question Text]
Answer in v_1 : [Answer] / Answer in v_2 : [Answer]

Reason:

- Why single-frame is insufficient: [Brief explanation]

7.3. Quality Control Process

We used 8 annotators for quality review (separate from the 10 video recorders and 4 human evaluators). Each instance was independently reviewed by 2 annotators, and video-question pairs in each instance were accepted only upon unanimous agreement. Instances were kept only if all constituent pairs passed. The inter-annotator agreement was 94.3%, and Cohen’s kappa was 0.835. The total review time was 62.2 hours (\sim 2-3 min per instance per annotator). We also report the per-category annotator agreement scores in Table 7; all three categories achieve substantial agreement ($\kappa > 0.8$). Event Attributes show the lowest agreement, consistent with the inherently more subjective nature of judging continuous properties such as speed and force, but still well above conventional thresholds for high-quality annotation. The final benchmark was then evaluated by 4 separate, independent

Table 7. Human annotation agreement across categories.

Category	Agreement	Cohen’s κ
Events	95.6%	0.88
Event Attributes	93.4%	0.84
Structural Event Logic	93.7%	0.81

annotators who achieved 98.2% I-Acc, confirming that it is high-quality and largely unambiguous.

7.4. Benchmark Statistics.

Data Statistics. This section reports dataset statistics for TimeBlind. TimeBlind includes both two-choice multiple-choice and binary questions in equal numbers. It is a short-video benchmark, with an average video duration of 8.49 seconds. We computed CLIP-based video similarity for all pairs by sampling frames at 1 FPS, extracting per-frame embeddings, and average-pooling within each video. The mean pair similarity is 0.971 (median 0.979), and 88.3% of instances have similarity above 0.95, confirming strong static consistency by design.

Table 8. Statistics of the TimeBlind benchmark.

Statistic	Value
Binary Question	1,200
Two-Choice Multiple-Choice Question	1,200
Avg. Video Length (sec)	8.49
Avg. Question Length (words)	20.54

Statistical significance. We report statistical significance tests below. Wilson 95% confidence intervals on I-Acc ($n = 600$): (1) Random chance: 6.3% [4.6%, 8.5%], (2) Best open-source (Molmo2-8B): 31.2% [27.6%, 35.0%], (3) Best overall (Gemini 3 Pro): 48.2% [44.2%, 52.2%], (4) Human: 98.2% [96.8%, 99.0%].

All intervals are non-overlapping, confirming that performance differences across tiers are statistically significant. Since all models are evaluated on the same 600 instances, we use the paired McNemar’s test for pairwise comparisons. The gap between Gemini 3 Pro (48.2%) and human performance (98.2%) is highly significant ($p = 6.69 \times 10^{-81}$). The gap between Gemini 3 Pro and the best open-source model, Molmo2-8B (31.2%), is likewise significant ($p = 1.67 \times 10^{-12}$). A one-sample exact binomial test of Gemini 3 Pro against random chance (6.3%) confirms performance significantly exceeds chance ($p = 1.91 \times 10^{-178}$). Even the small 1.9% gap between GPT-5 (46.3%) and Gemini 3 Pro (48.2%) correctly registers non-significant (McNemar’s $p = 0.441$), demonstrating that our sample size has appropriate discriminative power: it captures meaningful gaps while

not over-claiming marginal ones.

7.5. Ablation Study

FPS Ablation. In this section, we study how increasing the density of spatiotemporal sampling affects model performance. We evaluate the advanced closed-source models GPT-5 and Gemini 3 Pro and the open-source models Qwen3-VL-235B and Molmo2-8B at 1, 5, and 10 FPS. As shown in Table 9, denser sampling provides only modest gains. The best result (Gemini 3 Pro at 5 FPS, 56.2% I-Acc) still leaves a 42.0% I-Acc gap compared with human performance. Notably, performance does not monotonically improve—Qwen3-VL-235B peaks at 5 FPS and drops at 10 FPS, likely because longer input sequences degrade reasoning. These results confirm that temporal understanding deficits in TimeBlind are fundamental, rather than artifacts of low sampling rates.

Table 9. **FPS Ablation.** Performance under different frame sampling rates, with results reported in instance accuracy (I-Acc). The best result, Gemini 3 Pro at 5 FPS (56.2% I-Acc), still trails human performance by 42.0% I-Acc, and performance is not monotonic: Qwen3-VL-235B peaks at 5 FPS but declines at 10 FPS, likely due to degraded reasoning from longer input sequences.

Model	1 FPS	5 FPS	10 FPS
Human	98.2%	98.2%	98.2%
Qwen3-VL-235B	25.8%	35.0%	26.5%
Molmo2-8B	31.2%	31.7%	29.2%
GPT-5	46.3%	51.0%	49.3%
Gemini 3 Pro	48.2%	56.2%	49.6%

Video Similarity. We computed CLIP-based video similarity for all video pairs by sampling frames at 1 FPS, extracting frame-level embeddings, and average-pooling them within each video. We further evaluate GPT-5 on two similarity-based subsets, namely the top 30% highest-similarity pairs and the bottom 30% lowest-similarity pairs. As shown in Table 10, GPT-5’s performance is comparable across the highest- and lowest-similarity subsets (45.7% vs 47.8%), indicating that the benchmark’s difficulty does not depend strongly on residual visual differences between paired videos. Combined with the shortcut analysis in Table 4 (single-frame, shuffled-frame, and language-only baselines all near random chance), this confirms that TimeBlind is robust to subtle residual visual cues.

7.6. Taxonomy

In this section, we introduce the detailed taxonomy of TimeBlind as summarized in Table 11. The taxonomy is organized into three hierarchical levels: Event, Event Attribute, and

Table 10. **GPT-5 Performance on different similarity-based subsets.** The table reports GPT-5’s I-Acc results on the high-similarity subset, the low-similarity subset, and the full benchmark. GPT-5’s performance is comparable across the highest- and lowest-similarity subsets (45.7% vs 47.8%), indicating that the benchmark’s difficulty does not depend strongly on residual visual differences between paired videos.

Subset	I-Acc
Top 30% (highest similarity)	45.7%
Bottom 30% (lowest similarity)	47.8%
Full benchmark	46.3%

Structural Event Logic. For each level, we provide its formal definition, constituent sub-categories, and illustrative examples described through paired videos within the same instance.

7.7. Failure Cases


For each failure case in Figure 4–Figure 10, we present failure cases for state-of-the-art models, including closed-source models (GPT-5 and Gemini 3 Pro) and open-source models (Qwen3-VL-235B and Molmo2-8B). For each failure case, we show one question, two videos with their corresponding captions, and each model’s answer to the question for both videos. We indicate whether the answer is correct or incorrect using check marks and cross marks.


Table 11. **Hierarchical taxonomy of TimeBlind for video understanding.** The taxonomy is structured into three levels—Event, Event Attribute, and Structural Event Logic. For each level, we provide its conceptual definition, associated subcategories, and representative paired video descriptions.

Taxonomy	Definition	Sub-Categories	Example: Paired Video Descriptions
Event	The atomic units of temporal understanding, encompassing fine-grained action phases and diverse state transitions.	Fine-grained Action Subtle differences in action phase that are difficult to resolve from a single frame.	<ul style="list-style-type: none"> • <i>The man opens the door.</i> • <i>The man closes the door.</i>
		State Transitions Tracking attribute changes over time (e.g., color, size, shape).	<ul style="list-style-type: none"> • <i>The apple turns from red to yellow.</i> • <i>The apple turns from yellow to red.</i>
Event Attribute	The properties that characterize how an event unfolds over time.	Speed The rate at which an action is performed or an object moves over time.	<ul style="list-style-type: none"> • <i>The man quickly opens the door.</i> • <i>The man slowly opens the door.</i>
		Force The intensity of an action as reflected in visible physical effects (e.g., acceleration, impact strength, or deformation).	<ul style="list-style-type: none"> • <i>The man gently opens the door.</i> • <i>The man forcefully opens the door.</i>
		Magnitude The spatial extent or scale of motion or deformation.	<ul style="list-style-type: none"> • <i>The person waves slightly.</i> • <i>The person waves broadly.</i>
		Duration The length of time an event lasts.	<ul style="list-style-type: none"> • <i>The person holds the cup for a long time.</i> • <i>The person holds the cup briefly.</i>
		Direction The trajectory or spatial direction of motion.	<ul style="list-style-type: none"> • <i>The person moves toward the right side of the frame.</i> • <i>The person moves toward the left side of the frame.</i>
		Repetition The number of times an event occurs.	<ul style="list-style-type: none"> • <i>The person picks up the cup three times.</i> • <i>The person picks up the cup twice.</i>
Structural Event Logic	The higher-level rules governing how events relate to and compose with one another.	Temporal Topology Interval relations between event intervals (Allen’s interval algebra), such as before, meet, overlap, start, during, finish, and equal.	<ul style="list-style-type: none"> • <i>The person picks up the book and cup at the same time.</i> • <i>The person picks up the book and subsequently picks up the cup after a short delay.</i>
		Causal Contingency The causal dependencies between events.	<ul style="list-style-type: none"> • <i>The first press causes the light to turn off.</i> • <i>The second press causes the light to turn off.</i>
		Cross-Event Comparison The comparative reasoning over attributes of different events within the same video.	<ul style="list-style-type: none"> • <i>The person holds the cup longer than the book.</i> • <i>The person holds the cup for less time than the book.</i>

Examples 1

Question: How does the person's speed change in the video?
 A) Gradually accelerating. B) Gradually decelerating

Video 1: A girl walks in a park and gradually speeds up.


Video 2: A girl walks in a park and gradually slows down.
















 GPT-5:	Video 1  Video 2 	 Gemini-3-Pro:	Video 1  Video 2 
 Qwen3-VL-235B:	Video 1  Video 2 	 Molmo2-8B:	Video 1  Video 2 

Figure 4. Failure Case 1

Examples 2

Question: Does the basketball shot go in?
 A) The shot is made. B) The shot is missed.

Video 1: A person shoots a basketball, and it goes right into the hoop.


Video 2: A person shoots a basketball, but it does not enter the hoop.















 GPT-5:	Video 1  Video 2 	 Gemini-3-Pro:	Video 1  Video 2 
 Qwen3-VL-235B:	Video 1  Video 2 	 Molmo2-8B:	Video 1  Video 2 

Figure 5. Failure Case 2


Examples 3

Question:
Do the bottom lights turn on after a pause once the top lights go out?

Video 1: There is a small time gap between the top lights turning off and the bottom lights turning on.



Video 2: The bottom lights turn on as soon as the top lights turn off.








 GPT-5:	Video 1	✓	Video 2	✓	 Gemini-3-Pro:	Video 1	✓	Video 2	✓
	 Qwen3-VL-235B:	Video 1	✗	Video 2		✗	 Molmo2-8B:	Video 1	✓

Figure 6. Failure Case 3


Examples 4

Question:
Does the number of blue darts on the dartboard decrease?

Video 1: The number of blue darts on the dartboard increases.



Video 2: The number of blue darts on the dartboard decreases.








 GPT-5:	Video 1	✓	Video 2	✓	 Gemini-3-Pro:	Video 1	✓	Video 2	✓
	 Qwen3-VL-235B:	Video 1	✓	Video 2		✗	 Molmo2-8B:	Video 1	✗

Figure 7. Failure Case 4


Examples 5

Question: Are the lychees put into a plastic bag only once?

Video 1: Put the lychees into a plastic bag once.



Video 2: Put the lychees into the plastic bag three times.








 GPT-5:	Video 1	✗	Video 2	✓	 Gemini-3-Pro:	Video 1	✗	Video 2	✓
	 Qwen3-VL-235B:	Video 1	✗	Video 2		✓	 Molmo2-8B:	Video 1	✓

Figure 8. Failure Case 5


Examples 6

Question: Do the girl and the boy both pick up the same book?

Video 1: The girl and the boy both pick up the same book.



Video 2: The girl and the boy pick up different books.








 GPT-5:	Video 1	✓	Video 2	✓	 Gemini-3-Pro:	Video 1	✓	Video 2	✓
	 Qwen3-VL-235B:	Video 1	✗	Video 2		✓	 Molmo2-8B:	Video 1	✗

Figure 9. Failure Case 6


Examples 7

Question: Do the wiping strokes gradually expand?

Video 1: The wiper's strokes get bigger and bigger.



Video 2: The wiper's strokes get smaller and smaller.







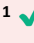







 GPT-5:	Video 1 	Video 2 	 Gemini-3-Pro:	Video 1 	Video 2 
 Qwen3-VL-235B:	Video 1 	Video 2 	 Molmo2-8B:	Video 1 	Video 2 

Figure 10. Failure Case 7